

Assignment - 3

Machine Learning.

Akshay Tambe (apt321@nyu.edu)
Snahil Singh (ss11381@nyu.edu).

Q1. $E(S) = - \sum_{i \in L} \frac{N_i}{N} \log_2 \frac{N_i}{N}$

$$IG(S) = E(S) - \sum_{v \in V} \frac{|S_v|}{|S|} E(S_v)$$

a) $A \rightarrow 4+, 5-$

$B \rightarrow 3+, 6-$

$$E(S_A) = - \left\{ \frac{4}{9} \log_2 \frac{4}{9} + \frac{5}{9} \log_2 \frac{5}{9} \right\}$$

$$E(S_B) = - \left\{ \frac{3}{9} \log_2 \frac{3}{9} + \frac{6}{9} \log_2 \frac{6}{9} \right\}$$

For the example with $(4+, 5-)$ i.e. A will have higher entropy than B because it's more balanced than B, or $N_+ / (N_+ + N_-)$ is closer to $1/2$ for A than B. A has less dissimilarity & less impurity than B.

b) $IG(X_1) = E(S) - \sum_{v \in V} \frac{|S_v|}{|S|} E(S_v)$

$$N(F) = 4, N(T) = 3$$

$$E(S) = - \left(\frac{3}{7} \log_2 \frac{3}{7} + \frac{4}{7} \log_2 \frac{4}{7} \right) = 0.52 + 0.46 = 0.98522$$

$$E(S_{X_1=T}) = - \left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \right) = -(-0.52832 + (-0.3899)) = 0.918295830$$

$$E(S_{X_1=F}) = - \left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) = - \log_2 \frac{2}{4} = 1$$

$$IG = 0.98522 - \left[\frac{3}{7} \times 0.91829 + \frac{4}{7} \times 1 \right] = 0.020244.$$

c) From question b,

$$\text{Entropy}(S) = 0.985228139$$

$$\text{Entropy}(S_{x_1=T}) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right) = 0.918$$

$$\text{Entropy}(S_{x_1=F}) = -\left(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4}\right) = 0.8113$$

$$\text{Entropy}(S_{x_2=T}) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right) = 0.918$$

$$\text{Entropy}(S_{x_2=F}) = -\left(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4}\right) = 0.8113$$

$$\text{Entropy}(S_{x_3=T}) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right) = 0.918$$

$$\text{Entropy}(S_{x_3=F}) = -\left(\frac{3}{4} \log_2 \frac{3}{4} + \frac{1}{4} \log_2 \frac{1}{4}\right) = 0.8113$$

$$\begin{aligned} \text{Information gain}(x_1) &= 0.985228139 - \left[\frac{3}{7} \times 0.918 + \frac{4}{7} \times 0.8113 \right] \\ &= 0.1281995 \end{aligned}$$

Similarly

$$\text{Information gain}(x_2) = 0.1281995 \&$$

$$\text{Information gain}(x_3) = 0.1281995$$

Comparing all the info gain for x_1, x_2, x_3 , we observe that all are equal

Hence, splitting on x_1

$$\begin{array}{c}
 x_1 \\
 \swarrow \quad \searrow \\
 F \qquad T
 \end{array}$$

	x_1	x_2	x_3	x
x^1	F	F	F	+
x^3	F	F	T	-
x^4	F	T	F	+
x^7	F	F	F	+

	x_1	x_2	x_3	x
x^2	T	F	T	+
x^5	T	T	F	-
x^6	T	T	T	-

Calculation of entropy &
Info-gain for left subtree

$$\text{Entropy}(S'_{x_2=T}) = -\left(\frac{1}{1} \log \frac{1}{1} + 0\right) = 0$$

$$\begin{aligned} \text{Entropy}(S'_{x_2=F}) &= -\left(\frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3}\right) \\ &= 0.918 \end{aligned}$$

$$\begin{aligned} \text{Info gain}(x_2) &= -\left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4}\right) \\ &= -\left(\frac{3}{4} \times 0 + \frac{1}{4} \times 0.918\right) \end{aligned}$$

$$\text{Info-gain}(x_2) = 0.5813$$

Calculation of entropy &
info-gain for right subtree

$$\text{Entropy}(S'_{x_2=T}) = -\left(0 + \frac{2}{2} \log \frac{2}{2}\right) = 0$$

$$\text{Entropy}(S'_{x_2=F}) = -\left(\frac{1}{1} \log \frac{1}{1} + 0\right) = 0$$

$$\begin{aligned} \text{Info-gain}(x_2) &= -\left(\frac{1}{3} \log \frac{1}{3} + \frac{2}{3} \log \frac{2}{3}\right) \\ &= -\left(\frac{1}{3} \times 0 + \frac{2}{3} \times 0\right) \end{aligned}$$

$$\text{Info-gain}(x_2) = 0.918$$

$$\text{Entropy}(S'_{x_3=T}) = -\left(0 + \frac{1}{1} \log \frac{1}{1}\right) = 0$$

$$\text{Entropy}(S'_{x_3=F}) = -\left(\frac{3}{3} \log \frac{3}{3} + 0\right) = 0$$

$$\text{Info-gain}(x_3) = -\left(\frac{3}{4} \log \frac{3}{4} + \frac{1}{4} \log \frac{1}{4}\right) - 0$$

$$\therefore \text{Info gain}(x_3) = 0.8113$$

$$\text{Here } IG(x_3) > IG(x_2)$$

hence splitting left subtree
on x_3

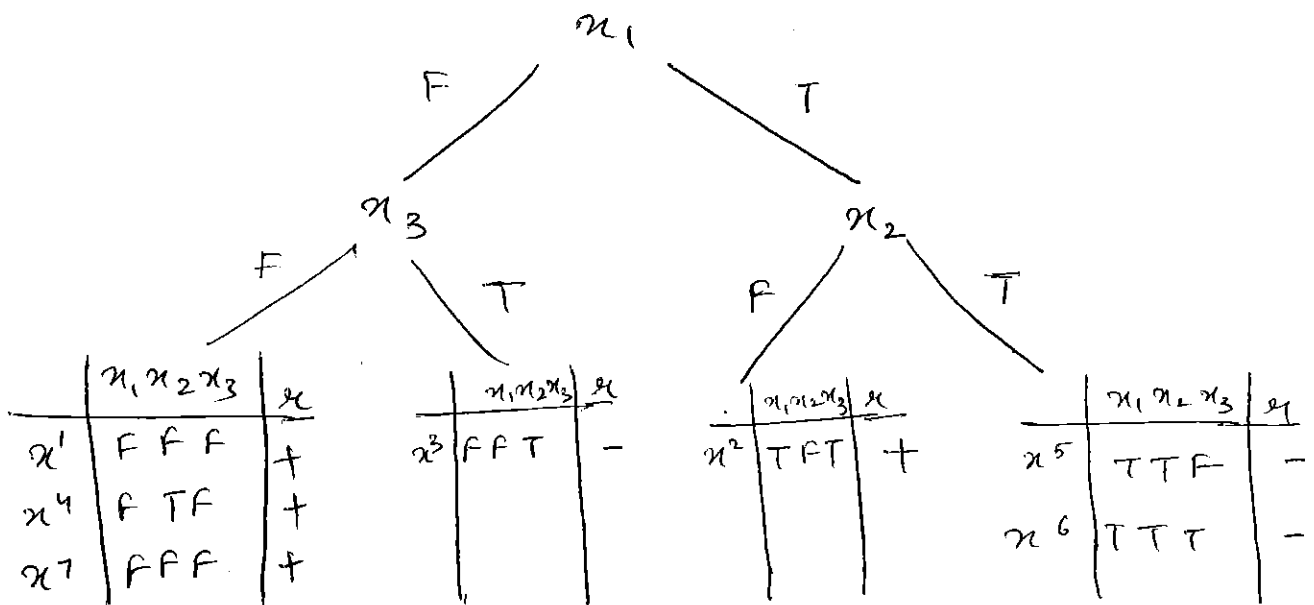
$$\text{Entropy}(S'_{x_3=T}) = -\left(\frac{1}{2} \log \frac{1}{2} + \log \frac{1}{2}\right) = 1$$

$$\text{Entropy}(S'_{x_3=F}) = -\left(0 + \frac{1}{1} \log \frac{1}{1}\right) = 0$$

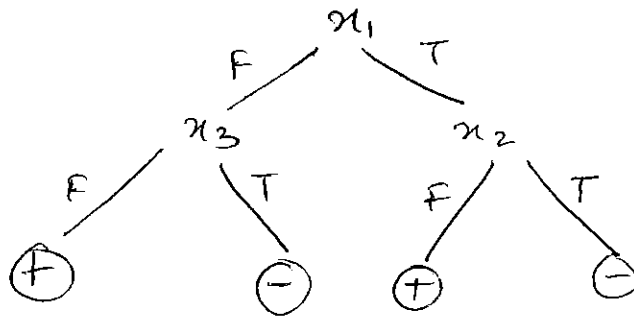
$$\text{Info-gain}(x_3) = 0.918 - \left[\frac{1}{3} \times 1 + \frac{2}{3} \times 0\right]$$

$$\text{Info gain}(x_3) = 0.58467$$

Here, $IG(x_2) > IG(x_3)$ hence
splitting right subtree on x_2



As entropy is 0 by splitting on x_3 on left subtree & x_2 on right subtree, we can clearly see from observations
Final tree -



d) $E H(X) = - \sum_{i=1}^n P[X=i] \log_2 P[X=i] \rightarrow$ expected no. of bits needed to encode

$$H(Y|X) = \sum_x P[X=x] * \left(\sum_y -P[Y=y|X=x] * \log_2 P[Y=y|X=x] \right)$$

M.I. b/w X and $Y = H(X) - H(Y|X) = H(X) - H(X|Y)$

$X \rightarrow$ value of attribute $x_i \in \{T, F\}$

$Y \rightarrow$ Label of example $\rightarrow +/-$

Need to compute $H(X)$ and $H(Y|X)$, $H(X) - H(Y|X)$

$$H(Y|X) = P(T) [-P(+|T) \log_2 P(+|T) - P(-|T) \log_2 P(-|T)] \\ + P(F) [-P(+|F) \log_2 P(+|F) - P(-|F) \log_2 P(-|F)]$$

$$= \frac{3}{7} \left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] + \frac{4}{7} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$H(Y|X) = -\frac{3}{7} \left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \right) - \frac{4}{7} \left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) \\ = -\frac{3}{7} (-0.92) = 0.396.$$

$$H(X) = - \left\{ \frac{4}{7} \log_2 \frac{4}{7} + \frac{3}{7} \log_2 \frac{3}{7} \right\} = 0.98.$$

$$H(X) - H(Y|X) = 0.584.$$

e) $Y \rightarrow 2$ possible labels.

$$\text{Entropy} = - \sum \frac{N_i}{N} \log \frac{N_i}{N}$$

Since all labels appear equally, we have $P(x) = 1/2$

$$= - \sum \frac{1}{2} \log \frac{1}{2} = - \log \frac{1}{2}$$